## NBER WORKING PAPER SERIES

# DOES RESEARCH SAVE LIVES? THE LOCAL SPILLOVERS OF BIOMEDICAL RESEARCH ON MORTALITY

Rebecca McKibbin Bruce A. Weinberg

Working Paper 29420 http://www.nber.org/papers/w29420

NATIONAL BUREAU OF ECONOMIC RESEARCH 1050 Massachusetts Avenue Cambridge, MA 02138 October 2021

Thanks to seminar participants at the Ohio State University and The University of Sydney. All errors are our own. McKibbin and Weinberg met through the NBER The Value of Medical Research network meetings, which were funded by the National Institute on Aging through Grant Number R24AG058049 and by the Economic and Social Research Council, through Grant Number ES/M008673/1 to the Institute for Fiscal Studies. This opportunity is gratefully acknowledged support from NIA, OBSSR, and NSF SciSIP through P01 AG039347; R01 GM140281, UL1 TR002733, NSF EHR DGE 1348691, 1535399, 1760544, 2100234; and the Ewing Marion Kauffman and Alfred P. Sloan Foundations. Weinberg was paid directly by NBER and his work at Ohio State was supported on a subcontract from NBER on P01 AG039347. The content is solely the responsibility of the author and does not represent the views of the NIH, IFS, or the National Bureau of Economic Research.

NBER working papers are circulated for discussion and comment purposes. They have not been peer-reviewed or been subject to the review by the NBER Board of Directors that accompanies official NBER publications.

© 2021 by Rebecca McKibbin and Bruce A. Weinberg. All rights reserved. Short sections of text, not to exceed two paragraphs, may be quoted without explicit permission provided that full credit, including © notice, is given to the source.

Does Research Save Lives? The Local Spillovers of Biomedical Research on Mortality Rebecca McKibbin and Bruce A. Weinberg NBER Working Paper No. 29420 October 2021 JEL No. 11,033,038

## ABSTRACT

This paper investigates the local impact of biomedical research on mortality in the USA. Causally estimating the marginal value of biomedical research is challenging due to a lack of micro data linking health outcomes to plausibly exogenous variation in research. We create a new linkage between a research database (PubMed) and administrative death records that enables research to be related to mortality at the geographic, disease and time level. We then estimate the marginal impact of biomedical research on mortality using hospital market (HRR) level shocks to research activity by disease. Our identification strategy builds on the literature on the dissemination of knowledge, specifically that of local knowledge spillovers. By utilizing variation across diseases, time and distance from research we control for additional trends relative to the current literature. Our results show that an additional research publication on average reduces local mortality from a disease by 0.35%. Our results also provide novel evidence that there are health benefits to the local communities (local spillovers) in which biomedical research is conducted.

Rebecca McKibbin Level 5, Social Sciences Build The University of Sydney Camperdown Camperdown, NSW 2050 Australia rebecca.mckibbin@sydney.edu.au

Bruce A. Weinberg The Ohio State University Department of Economics 410 Arps Hall 1945 North High Street Columbus, OH 43210 and NBER weinberg.27@osu.edu

# 1 Introduction

Despite large investments in biomedical research (the NIH budget alone is over \$40bn per year), there is surprisingly little causal research linking improvements in health to biomedical research. Over the past century the expansion of scientific understanding of disease has led to large changes in the practice of medicine and available treatments for many diseases. Murphy and Topel (2006) estimate that the contribution of medical research to reductions in mortality over the twentieth century are worth \$3.2 trillion per year. However it is difficult to evaluate the success of biomedical research on the margin in improving health outcomes because the results of scientific studies tend to be narrowly focused and incremental, and depend on the adoption of the results, creating a measurement and identification challenge. Most research on the impact of research on health is based on intermediate outcomes such as publications produced, patents registered or completed clinical trials. In this paper we link health (measured by mortality) to research output, specifically publications, and grant funding (an input) to evaluate the average incremental impact of biomedical research on its ultimate goal: improving health.

The empirical strategy we use in this paper leverages variation in the amount of biomedical research produced in different geographic regions, over time and across disease categories. This strategy is built on two related literatures: the production of science and the theory of technology diffusion. A core principle of the canonical model of technology diffusion (Rogers, 1962) is that that the adoption of new ideas does not occur simultaneously, rather, ideas diffuse gradually through social systems, which have a geographical component. This has been shown to hold empirically across a range of settings (Feldman and Kogler, 2010). In our setting, research and the application of research into practice are closely intertwined with academic medical centers that provide clinical care to patients playing a large role in the production of research and medical training. These multiple roles make them an ideal setting for facilitating local knowledge spillovers. If local spillovers exist we would expect to see a temporary health benefit to patients from being treated by physicians close to where research on their ailment is conducted because their treatment is impacted by the new knowledge earlier than those located further away. Agha and Molitor (2018) find evidence that physicians geographically proximate to pivotal trials are earlier adopters of new cancer drugs, lending support to our empirical strategy.

Estimating the relationship between local mortality and local research is challenging because it is likely dynamic — research spreads slowly over time — and potentially endogenous — researchers choose what to research, and they may consider changes in local population health in selecting their research agendas.<sup>1</sup> To capture these dynamics we utilize the Jordà (2005) local projection method to generate an impulse response for local mortality to a shock to local research.<sup>2</sup> Shocks to research

<sup>&</sup>lt;sup>1</sup>Their choice of research topic will be influenced by the availability of resources, such as grants, and local health needs as well as their area of expertise, which is a product of interests, training, and previous research activities.

 $<sup>^{2}</sup>$ This approach to estimating impulse responses has a number of advantages over traditional IRFs estimated using VAR models, notably, that they have more flexible functional forms, are robust to lag length mis-specification and can be estimated using an instrumental variables approach to shock identification.

activity are identified using instrumental variables. The instrument utilizes the differential impact of national shocks to research funding across diseases that impacts locations differently. The results show that a one percent increase in publications on a disease in a community leads to a cumulative 0.35% reduction in the local mortality rate for that disease in the local area where the research was conducted.

In addition to the identification challenges, our analysis poses a technical challenge — linking biomedical research to health outcomes. We create a linkage between ICD-10 codes, which are used to measure causes of death (our measure of health outcomes) and Medical Subject Headings (MeSH) codes, which are used to index the subject matter of biomedical research publications. This is done by using the NLM MeSH generator to convert text descriptions of the ICD-10CM codes into MeSH codes. Since grant data are not comprehensively indexed to topics over our long time horizon, we link the grant data to the publication output reported to the NIH and index the grants using the MeSH codes associated with the publications, with funding being allocated in proportion to the frequency of the MeSH code across all output produced from the grant. Our final dataset contains per capita mortality by hospital referral region for thirty-eight of the most common causes of death and the corresponding publications and grants that are related to these causes of death.

Estimating local spillovers of biomedical research to health is of interest in two literatures. First, it provides an estimate of the marginal value of biomedical research. Murphy and Topel (2006) estimate the value of the reduction in mortality over the twentieth century using a structural approach. This approach enables the potential benefits of biomedical research, however, they do not include estimates of the causal mechanism between investment in biomedical research and improvements on health, which requires micro-level data. Lichtenberg (2018a) examined this question in the setting of cancer and found that biomedical research did improve cancer survival. Our approach incorporates an additional dimension of variation — across geographic locations — a broader range of diseases and a different estimation strategy.

Our approach enables us to overcome a significant challenge faced by the literature – identifying the impact of biomedical research on health itself as opposed to intermediate steps in the translation process. Much of the empirical literature has focused on identifying the impact of grant funding on intermediate outcomes such as the number of patents (Azoulay *et al.*, 2018; Toole, 2007) or the number of pharmaceutical products approved (Sampat and Lichtenberg, 2011). Or the impact of intermediate products such as pharmaceuticals on health (Lichtenberg 2019a; 2019b; 2018b; 2013) However, beyond commercial technology like pharmaceuticals, and outside of cancer, there has been limited research on the impact of research on health. While these intermediate outputs are important, bridging the gap in evidence on the actual effect of biomedical medical research on health is critical in for directing public funding in biomedical research.

Second, this paper provides empirical evidence for knowledge spillovers in the health care setting. The literature on knowledge spillovers is, by now, large and the findings have implications as wide ranging as economic growth, urban agglomerations, and international trade (Romer, 1986; Lucas Jr, 1988; Glaeser *et al.*, 1992; Krugman, 1991). Perhaps the closest literature focuses on knowledge spillovers among researchers (Waldinger, 2010, 2012; Borjas and Doran, 2012, 2015; Ham and Weinberg, 2021) with Zucker *et al.* (1998) and Azoulay *et al.* (2010) focusing specifically on spillovers among and from biomedical researchers. Previous research has also estimates a range of benefits from universities as engines of innovation and growth for local communities (Bania *et al.*, 1993; Beeson and Montgomery, 1993; Saxenian, 1996; Moretti, 2004; Kantor and Whalley, 2009; Zolas *et al.*, 2015).

The existence of spillovers shows that there are unrealized gains from medical knowledge. With scientific advances appearing to be increasingly difficult (Bloom *et al.*, 2020), identifying potentially unrealized health gains from existing knowledge is an important opportunity to increase the returns on investment in science and improve population health. Our results indicate that there are consequential gains to be had in improving the dissemination of research findings among physicians.

# 2 Background

Biomedical research is a branch of science that studies the physiology and treatment of disease and illness. The umbrella of biomedical research includes a continuum of subjects that range from basic to translational to applied clinical research. Our analysis focuses on research that is likely to fall broadly into the category of clinical research, since this is the type of research that is most likely to influence clinical practice. Clinical research is conducted by scientists and clinicians in commercial and academic settings. While commercial research tends to focus on patentable discoveries, such as drug and medical device development, academic researchers and clinicians also conduct research that may have little or no commercial value, but may yield important changes in health care.

It is important to consider what drives the production and dissemination of research and how this might affect the identification of causal effects. Academic researchers have a great deal of autonomy over the the topics that they research. A large literature has examined what motivates their work. A broad categorization of the key factors are intrinsic motivation (or interest in the topic), reputational rewards and financial rewards (for example in the form of promotions)(Lam, 2011). However, whatever the motivation, producing any research requires (often substantial) resources. This includes (potentially) expensive equipment, salaries for staffing the laboratory with technicians, students and junior researchers (Stephan, 2010) and overhead to the university to cover the university's costs of maintaining facilities.

Academic biomedical research relies primarily on external funding, that is, other than a startup fund for newly hired faculty, research labs need to be financially viable independent from the university. This means that funding sources may play a role in determining research activity. In the US, the largest funder of biomedical science is the NIH. In FY2020, the NIH budget was \$41.7B, over 80% of which went to researchers at universities or other institutions outside of NIH.<sup>3</sup> There

<sup>&</sup>lt;sup>3</sup>Source: https://www.nih.gov/about-nih/what-we-do/budget. Accessed 3/10/2020

are two broad avenues of funding for extramural research: "requests for applications" (RFAs) and "investigator-initiated grants". The former are targeted at specific NIH priority topics while for the latter the scientist sets the subject. Projects are chosen for funding using a peer-review process. The process takes into account the potential of the project as well as the research history of the investigators (a measure of their expertise and potential to complete the project).

While NIH funding does respond to health needs – for instance, the GAO finds that there is a correlation between cause of death and funding – the response is not immediate nor is it deterministic. Specifically, the NIH sets research priorities each year in consultation with its 27 institutes and centers (ICs). A budget is then submitted to the Department of Health and Human Services (HHS). There is then a back and forth between the NIH and HHS before the budget is finally submitted to the Office of Management and Budget (OMB) and then to Congress for approval, which adds additional uncertainty. The time from priority setting to funding appropriations is typically 18 months. This long lag results in sticky and somewhat unpredictable funding for specific diseases.

Although funding is critical for the production of research, so too is expertise and knowledge. Myers (2020) and Azoulay *et al.* (2010) both examine how malleable research topics are to funding opportunities. Researchers need to have expertise and build up knowledge in order to advance science, so there is a lot of persistence in the topics researchers pursue. This persistence is important for our identification strategy, which relies on the allocation of shocks to funding for a particular disease being influenced to some extent by initial differences in research capacity on that disease.<sup>4</sup>

Commercial motivation likely also influences how research makes its way into practice. In the case of pharmaceuticals, the role of detailing (marketing directed at physicians) has been studied extensively (Ching and Ishihara (2010); Ching *et al.* (2016); Chintagunta *et al.* (2012)). However, when there is no financial incentive to advertise, knowledge likely spreads more slowly and informally. For example, publications that aggregate research results, clinical practice software or conferences could facilitate the spread of new ideas among clinical physicians. It is also reasonable to expect that some, though not all, doctors keep abreast of the latest literature. From a theoretical perspective, we would expect that knowledge and innovation spread within interpersonal networks (Coleman *et al.*, 1957) and geographic regions (Baicker and Chandra, 2010). Agha and Molitor (2018) investigate this in the context of new cancer drugs. Their results show that patients treated in the region where the first author of a study is located are substantially more likely to receive treatment with a new drug within the first two years following a drug's FDA approval but they do not find significant impacts on mortality. One reason might be that cancer drugs appear to have very low efficacy profiles, making it difficult to detect effects in observational data.

<sup>&</sup>lt;sup>4</sup>Sattari and Weinberg (2017) identify the NIH IC that provide the most funding to each researcher in his or her first year of NIH funding and show that on average 80% of funding for researchers comes from the IC that provided the most funding in their first year of NIH support. Even 30 years after initially receiving NIH funding over 60% of funding comes from the IC that provide the most funding in the researcher's first year.

# 3 Data

Data from administrative death records is linked with data on publications and NIH grant funding using a newly created cross-walk that maps MeSH codes with ICD-10 codes and indexes grants with MeSH codes. The dataset contains information on publications, grants and deaths by cause of death at the hospital market level, disease and year.

### 3.1 Measuring Mortality

Mortality is measured using data from the National Center for Health Statistics Detailed Multiple Cause of Death (MCOD) Research Files (1999-2017), which are administrative data drawn from death certificates. Our primary mortality outcome is age-adjusted years of potential life lost per capita. Years of potential life lost (YPLL) is a measure of premature mortality that places greater weight on deaths at younger ages relative to older ages. We use this commonly used measure for two reasons. Firstly, it captures the fact that medical advances could greatly improve outcomes for patients with a disease, even though they will still die with or of the disease. Secondly, it reflects the fact that everyone will inevitable die from some cause and when it comes to the allocation of resources towards preventing these causes, there is a societal preference for finding ways to prevent death those who have longer remaining life expectancies.

Potential life lost is computed for each decedent by subtracting their life expectancy at the age at which they died, given the year in which they were born. Life expectancies were obtained from the United States Mortality DataBase (1999-2017) by individual age. We age-adjust the mortality outcomes to account for changes in the age profile of the local population over time. This is done by re-weighting the mortality measures using a fixed population distribution. We use the 2000 U.S. Standard Population, 19 age groups.<sup>5</sup> The age-adjusted mortality is then converted to a per capita measure using county level population data from the SEER database.<sup>6</sup> Our unit of geographic analysis is a hospital referral region (HRR), which is a commonly used geographic definition in health services and health economic research (Kibria *et al.*, 2013). The mortality outcomes are aggregated from the county level to the HRR level using the cross-walk provided by the Dartmouth Atlas and the county of residence at death.<sup>78</sup>

The analysis requires the mortality outcome to be measured by disease category, and disease category to be linkable with research. In order to ensure the disease categories are relevant to deaths, we begin by using the NCHS 113 most common causes of death, which are groupings of ICD10-CM codes. We exclude causes of death due to injury, residual categories that would be very difficult

<sup>&</sup>lt;sup>5</sup>Source: https://seer.cancer.gov/seerstat/tutorials/aarates/step1.html.

<sup>&</sup>lt;sup>6</sup>Source: https://seer.cancer.gov/popdata/download.html.

<sup>&</sup>lt;sup>7</sup>Source: https://data.dartmouthatlas.org/supplemental/.

 $<sup>^{8}</sup>$ We check that our results are robust to using the county of occurrence. In cases where a county does not map to a single HRR, we count the deaths in both HRRs.

to match and categories where the mapping MeSH and ICD10-CM is not unique.<sup>9</sup> Further details are provided in Appendix A. Our analysis is based on thirty-eight categories of diseases, which are listed in Table A.1.

Deaths are classified into these categories using the ICD-10CM codes listed in the multiple cause of death fields on the death certificate. Death certificates contain both a single underlying cause of death as well as up to twenty factors that contributed to the death. We count a death towards every category that was a contributing factor. This enables us to capture changes in premature mortality for diseases that may not commonly be direct underlying causes but still contribute to the deterioration of health.

The first column of Table 1 shows descriptive statistics of the primary outcome variable - years of potential life lost per capita.

## 3.2 Measuring Local Research Activity

The key explanatory variable for our analysis is a measure of research activity for each year, disease and HRR (geographic location). We create two measures of this for each disease, HRR and year: a measure of research outputs based on publication data and a measure of research inputs based on receipt of grant funding from the NIH.

Measures of publications are created using the NIH National Library of Medicine (NLM) PubMed database, which is an index of over 32 million biomedical research publications. We require the measures to made by geographic location, disease and time. We attribute the paper to an HRR using the first affiliation of the first author and to a year based on publication year.<sup>10</sup> To categorize local publications by disease we utilized the Medical Subject Headings (MeSH) attached to the publications. MeSH codes are a controlled and hierarchically organized vocabulary that is used to index the subject matter of biomedical information. Publications catalogued in PubMed are tagged by independent, professional coders at the NLM with one or more MeSH codes. We create a cross-walk between the disease category groups based on the ICD10 codes used in the mortality data with MeSH headings. This is done by taking the text attached to the three digit ICD-10 codes, which comprise a disease grouping in our mortality data, and indexing this with MeSH terms using the "MeSH on demand" tool provided by the NLM. We use terms that fall under the level 1 heading "Diseases". Following Packalen and Bhattacharya (2017) this should result in a measure of disease-based research that is a proxy for clinical research. We include all terms below the returned terms in the MeSH hierarchy and manually verify the relevance of the selected MeSH terms. This process and a cross-walk of MeSH terms to ICD10 causes of death is provided in the Appendix. Summary statistics for the publication measure are provided in Table 1.

 $<sup>^{9}</sup>$ We also exclude three disease categories that are relatively rare causes of death and have very small, and hence concentrated research locations. The results are robust to the inclusion of these diseases (see Appendix C.2).

<sup>&</sup>lt;sup>10</sup>In biomedical research the first author is typically the researcher who contributed the most work to the publication while the last author is typically the senior PI. We use the first author because prior to 2014, only the affiliation of the first author is provided in PubMED.

	$\mathbf{YP}$ (yea	LLrs)	Publica (No	ations .)	<b>Fund</b> (00	$\log^a ($
	Mean	SD	Mean	SD	Mean	SD
Cancer						
Bladder cancer	8.3	8.2	1.9	4.9	53	170
Brain cancer	14.4	15.3	7.6	17.2	391	1,013
Breast cancer	33.2	32.1	17.3	37.0	1,042	2,330
Cervical cancer	4.9	6.3	2.4	5.4	98	292
Colorectal cancer	38.2	38.1	8.2	17.5	520	1,208
Female gyn cancer	5.2	5.5	2.1	5.3	59	187
Kidney cancer	10.4	11.3	3.2	8.0	61	184
Leukemia	17.9	17.8	6.2	14.7	310	742
Liver cancer	13.0	13.5	4.3	10.4	159	385
Lung cancer	110.0	110.1	8.6	19.2	435	1,046
Melanoma	7.3	7.9	4.3	10.5	198	512
Mouth cancer	7.5	8.3	2.2	5.3	67	205
Multiple myeloma	7.5	7.6	1.8	5.4	60	198
Non-hodgkin lymphoma	15.3	15.4	4.0	9.4	158	457
Oesophagus cancer	10.4	10.5	1.5	3.8	37	147
Ovarian cancer	10.0	10.1	3.7	8.5	193	474
Pancreas cancer	21.8	21.0	3.7	8.8	187	471
Prostate cancer	17.5	17.8	10.4	23.3	542	1.254
Stomach cancer	7.4	8.1	1.2	3.2	28	109
Other Diseases						
Acute myocardial infarction	122.1	144.9	5.9	13.2	287	749
Alcoholic liver disease	21.4	33.3	0.5	1.5	46	154
Alzheimer	33.6	35.9	6.7	15.2	724	1,921
Anemias	25.4	27.4	4.9	10.3	248	690
Atherosclerosis	20.9	26.4	2.3	6.8	286	806
Cerebrovascular diseases	117.9	119.6	13.8	28.4	734	1,695
Complications of medical and surgical care	23.5	26.2	14.1	30.6	123	336
Diabetes	143.5	147.5	17.7	34.1	1,281	2,931
HIV	11.1	16.2	17.4	43.5	2,735	7,732
Intestinal infectious diseases	4.9	6.6	5.0	9.5	278	720
Malnutrition	13.9	16.4	2.9	6.0	149	399
Meningitis	2.0	3.7	1.1	2.5	54	204
Nephritis	115.5	116.2	4.8	9.8	178	473
Other heart disease	415.4	426.9	8.5	18.0	373	940
Parkinson disease	13.1	13.1	3.4	7.3	263	695
Peptic ulcer	4.3	5.4	0.5	1.3	4	21
Respiratory diseases	156.7	161.9	10.9	21.7	787	2,104
Septicemia	93.9	94.9	4.0	8.0	174	416
Viral Hepatitis	14.9	18.1	4.6	10.5	302	911
Total	45.1	117.3	5.9	17.2	358	1,669
Sample size <sup>b</sup>	220.932					

Notes: (a) Funding is converted into real dollars 2000 using NIH's Biomedical Research and Development Price Index (BRDPI). (b) Sample comprises 309 Hospital Referral Regions, 19 years and 38 diseases(T=19, N=11,934 (HRRs:306 x 10.0 m)) and the same set of the same set

Diseases: 38))

#### Figure 1: Research Activity



Notes: These maps show the geographical share of publications produced and NIH funding received in 1999, which is the first year of our analysis period. The color groupings represent quintiles. Publications are attributed to an HRR based on the affiliation of the first author. Grants funding is first allocated among publications that are linked to the grant and from there assigned to MeSH categories and geographic location in the same was a publications. Since grants can fund multiple publications and funding is disbursed over multiple years, the funding received in each year is assigned in equal share to all publications that have not yet been published at the time.

Data on NIH grants is taken from NIH's ExPORTER tool. NIH's ExPORTER tool provides back end access to its Research Portfolio Online Reporting Tools Expenditures and Reports (Re-PORTER) database. This dataset contains detailed information on all research projects and subprojects funded by NIH since 1985. Here we focus on support for extramural research, including research grants (activity codes beginning with "R"), cooperative agreements (U01) and early career awards (K99)). More than 83% of the projects listed in the ExPORTER data are research grants.

The grants data are not indexed by topic by the NIH until 2008. To create disease groupings for grants we merge the grants data with the publication data based on the publication id's reported to the NIH by the grant awardee. For each grant, we attribute funding equally to each publication that had not yet been published in the year the funding was disbursed and is published within 10 years of the funding disbursement.<sup>11</sup> The funding is then categorized into the HRR-disease level following the approach used to categorize the publications themselves. However, the year of the funding is the financial year in which the funding was disbursed. Funding is converted into real dollars 2000 using NIH's Biomedical Research and Development Price Index (BRDPI). Summary statistics by disease are provided in Table 1.

Figure 1 shows quintiles of the share of research publications and funding by hospital market region in 1999, the first year of our outcome variable. Research and funding a distributed across the US, with the largest shares in New York City, New Jersey, Boston and Los Angeles. Research and funding are highly correlated (0.94).

One source of variation in the analysis is variation across diseases within HRRs. Figure 2(a) shows the coefficient of variation in funding across diseases within an HRR. This is computed by first

<sup>&</sup>lt;sup>11</sup>For robustness we remove the ten year restriction, this makes little difference to the results.

#### Figure 2: Coefficient of Variation of NIH Funding



Notes: These maps shows quintiles of the coefficient of variation in the funding data. Figure (a) shows the coefficient of variation within each HRR across diseases. It is computed using the average funding received in each location for the disease over the time period 1999-2017. Figure (b) shows the coefficient of variation within each HRR across time. It is computed using the average share of funding across diseases within the HRR in each year. The coefficient of variation is then computed using this HRR-time variable.

averaging the share of disease funding in each year that an HRR receives across the time period 1999-2017. The mean and standard deviation used to compute the the coefficient of variation are then computed for each HRR. Hence they show variation across diseases within an HRR for the average funding share over time. Comparing this with Figure 1, it appears the variation across diseases within an HRR is larger in HRRs with smaller research programs. Figure 2(b) shows the variation within an HRR over time. This is computed by first taking the average funding share received by each HRR for each disease in each year across diseases. The coefficient of variation is then computed using the mean and standard deviation for each HRR (hence the variation is across time). Again, much of the variation over time is in HRRs with lower shares of research inputs and outputs in 1999.

## 4 Empirical Strategy

Our empirical strategy is based on the idea that physicians who are geographically proximate to the developer of an idea are more likely to be early adopters of that idea. While over time, good ideas spread through the dissemination of research findings (e.g via publications and professional networks), earlier access to medical advances resulting from geographic proximity to research that targets a particular health disease yields a health advantage to people with that disease in the locations where research is conducted. This short-run effect is what we seek to estimate.

In an ideal experiment, each year early access to research results on each disease would be randomly allocated to a geographic region. The research would then be unrelated to unobservable characteristics of the patient population in a particular region both present and past. We could then estimate the impact of this earlier access to research results relative to places that did not receive earlier access. In our non-experimental setting, research results appear first in the locations where the local researchers made the choice to undertake that particular research. This means that research output could be correlated with unobservable characteristics of the location in relation to the particular disease. For instance if there is a particular health problem in a location, it might be more salient to researchers in that area and may pique their interest. It is also potentially easier to conduct clinical research when there is a large patient pool available to study.

To address this issue we use an instrumental variable approach. We apply two distinct but related instruments that seek to emulate an unexpected windfall in research funding, and hence research output. The first uses the differential impact of national shocks to funding allocations for diseases groups, on locations based on their estimated ability to capture funding on the margin. The second, uses a natural experiment that resulted in windfall funding. In 2009, the NIH received a large shock to its budget from the American Recovery and Investment ACT. A substantial share of this funding was used to extend the payline on grant applications received in 2008 and 2009. We use this shock to support our main analysis, as the follow-up period is short.

#### 4.1 Estimation using local projection

Our goal is to estimate the relationship between local mortality and local research. The key features of this relationship are that it is dynamic and likely endogenous. Using the Jordà (2005) local projection method, we estimate an impulse response function (IRF) of local mortality to a shock to local research activity. This method has several advantages over a traditional VAR specification. Notably, it is more flexible and robust to misspecification of lags, moreover it does not constrain the shape of the IRF and it easily accommodates an instrumental variables approach.

Equation 1 shows the estimating equations. We estimate Equation 1 at each h step ahead to generate the impulse response function of local mortality (YPLL) to a one percent shock to local research.

$$log(m_{ld,t+h}) = \alpha^{h} + \beta^{h} log(R_{dl,t-1}) + \psi^{h}(L) z_{t-1} + \delta^{h}_{lt} + \gamma^{h}_{dt} + \eta^{h}_{ld} + \epsilon_{ld,t+h}, for \ h = (0, 1, 2, \dots 10) \ (1)$$

Age adjusted years of potential life lost per capita  $(m_{ldt})$  in a location l for a disease d at time t+h are a function of the shock to research R conducted in that location on that disease in at time t-1.  $\psi^h(L)$  is a polynomial of the lag operator,  $z_{t-1}$  is a vector of control variables, which includes additional lags of research. The model includes fixed effects for location\*time  $(\delta_{lt})$ , which account for changing demographics in a local area; disease\*time  $(\gamma_{dt})$ , which captures national trends in the disease; location\*disease  $(\eta_{ld})$ , which capture level differences in the mortality rate for a particular disease in a particular place. Put differently, the unit of analysis is a disease\*time\*location and our model includes all three pairs of two-way interactions between location, time, and disease. The coefficients of interest are the  $\beta^h$ , which are the effect of a 1% increase in local research activity on

the local mortality rate relative to the mortality rate for the same disease in other locations and relative to the same location and other diseases h periods after the shock.

Figure 3 shows the estimates of the  $\beta^h$  for a one percent shock to research activity (measured in the year prior to the year mortality is measured), for h = (0, 1, 2, ...10).<sup>12</sup> Three measures of publications are included: all publications, publications that acknowledge NIH funding, publications in journals with an Scimago Journal Rankings (SJR) in the top quartile as a measure of publication quality. The fourth figure shows the response to a one percent increase in local NIH funding received for a disease.

The results for "all publications", which is our main measure of research output, show that local research publications for a particular disease (if anything) are associated with higher local mortality from that disease. The relationship is negative for research that is NIH funded and for research published in a more highly rated journal than for the more general measure of any publication. The relationship between NIH funding and mortality is negative and smaller than for NIH publications. This reflects that a one percent increase in funding results is a less than one percent increase in output.

#### 4.2 Identifying a shock to research

The previous results are identified under the assumption that after controlling for fixed effects the remaining variation in publications is exogenous to local health diseases past and present. This is naturally a strong assumption given that research topic is a choice made by researchers and they may be influenced by local health, particularly in the past given the timing of funding decisions relative to the health outcome being observed. To address this concern we propose two related but distinct instruments. The first uses the differential impact of national shocks to funding allocations for disease groups, on locations based on their estimated propensity to capture funding on the margin. The second uses a large temporary increase in the NIH budget – the American Recovery and Reinvestment ACT (2009) as shock to funding. The former has the advantage of enabling the study of a longer period, the latter is a more distinct shock, but has a shorter follow-up period.

## 4.2.1 Capturing NIH funding on the margin

The idea of the first instrument is that locations have differential exposure to national funding shocks because of pre-existing relationships and research backgrounds that make them more effective at obtaining increases in NIH funding on a specific disease. This means that when the NIH has some additional funding for a specific disease, some areas have a higher chance of capturing this funding than other areas for reasons unrelated to the local health of the population. This could occur, for example, if there is a Matthew Effect (Merton, 1968).

 $<sup>^{12}\</sup>mathrm{Full}$  regression results for the outcome "all publications" are shown in Table A.2.



Figure 3: Response of Age-adjusted YPLL per capita to a 1% Shock to Research

Notes: This figure shows the estimates of  $\beta^h$  in Equation 1 and the 95% confidence interval. These results are generated by a separate regression for each h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the dependent variable of interest is the inverse sine of research (publications or funding) in a location (hospital market) related to each disease in each year. Regressions control for all three pairs of fixed effects and ten lags of publications. Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

It is not anticipated that this propensity to capture funding is randomly distributed across locations — to a large extent the location of research depends on the location of universities and academic medical centers. Instead, we will estimate the extent to which shocks to aggregate funding for a disease result in greater mortality reductions for a disease in places that are expected to receive a larger share of the funding increase relative to places that are predicted to receive less of the shock. This strategy relies on the exogeniety of the national funding shocks and an estimated relationship between research, funding and mortality.

This instrument is constructed through estimation. As shown in Equation 2, we estimate location specific elasticities of funding for research on a disease that give a location's propensity to capture an extra dollar of NIH funding per capita allocated to individual diseases. We estimate this "capture rate" by regressing the inverse sine (interpreted in the same way as a log) of funding in location lon disease d at time  $t - log(funding_{ldt})$  — on the inverse sine of national funding on that disease at the that time,  $log(funding_{dt})$ . Both the local and national funding are specified in per capita terms. In this formulation,  $\beta$  is the elasticity of local funding to national funding. The idea is to isolate variation in local funding for a disease that come from fluctuations in national funding that are plausibly exogenous to local health shocks. Our estimates identify a local average treatment effect – it is plausible that funding is allocated on the basis of the promise of the research project and hence additional funding to the same topic or scientists may have a diminishing return.

Equation 2 shows the specification used to estimate the elasticities for the in the instrument  $(\hat{\beta}_{dl})$ 

$$log(funding_{ldt}) = \beta_{dl}log(funding_{dt}) + \kappa_d log(death_{ldt}) + \eta_{ld}log(death_{dt}) + \alpha_{dl} + \alpha_{lt} + \varepsilon_{ldt}$$
(2)

The instrument (Equation 3) is then constructed as the product between the propensity to capture funding  $(\hat{\beta}_{dl})$  and national funding for that disease. The national funding for a disease  $(\frac{funding_{dt}}{funding_t} \times \overline{funding_t})$  is computed as the disease's share of national funding in each period, holding total national funding constant at the mean over the period 1999-2017. This removes variation from growth or real declines in the NIH budget over time.

$$funding\_inst_{dlt} = \hat{\beta}_{dl}log(\frac{funding_{dt}}{funding_t} \times \overline{funding_t})$$
(3)

**Exclusion Restriction** The exclusion restriction requires that the unexplained components of the local mortality rate are uncorrelated with the product of national funding shocks and the estimated skill at capturing this funding. This means that locations that have a higher funding elasticity for a particular disease should not face a systematically different health response to national funding shocks to that disease. Put differently, national funding for diseases should not be chosen strate-gically based on local health trends that relate to the local propensity to capture national funding after controlling for our three, two-way fixed effects. Intuitively it seems plausible that NIH funding variations would be be based on the perceived quality of the science being produced, the national disease burden, and idiosyncratic bureaucratic factors.

#### Figure 4: Local Disease Funding Shares



Notes: This figure shows the distribution of actual funding shares captured for each disease-location-year.

To provide support for the exclusion restriction we show two descriptive statistics. First, that no region captures a large share of NIH funding for a disease, making their individual health diseases unlikely to carry significant weight in budgeting decisions at the NIH. This is shown in Figure 4. The largest funding share captured by any region for any disease is 28%. However the 99th percentile sits at 4.47%. Second, we estimate the following regression

$$log(\frac{funding_{dt}}{funding_t}) = \beta_k log(ms_{d,t+k}) + \delta_t + \delta_d + \epsilon_{dt} \qquad k \in [-3,3]$$

$$\tag{4}$$

where  $log(\frac{funding_{dt}}{funding_t})$  is the share of national funding allocated to disease d at time t (this mechanically sums to one in each period) and  $ms_{d,t+k}$  is the share of age adjusted years of potential life lost per capita for a disease d at time t + k. We include fixed effects for time  $(\delta_t)$  and disease  $(\delta_d)$ . The purpose of this regression is to examine whether the funding share allocated to a disease is related to fluctuations in mortality after controlling for fixed effects. Figure 5 shows the results. The average national shock at the disease level is not explained by the past three years or future three years of average mortality for that disease.

**First Stage** The first stage is estimated as follows:

$$log(R_{dl,t}) = \alpha + \beta funding\_iv_{dl,t-k} + \psi(L)z_{t-1} + \delta_{lt} + \gamma_{dt} + \eta_{ld} + \epsilon_{ld,t}$$
(5)

where  $R_{dl,t}$  is the research conducted in location (l) on a disease (d) at time t and funding\_iv\_{dl,t-k} is the instrument described in Equation 3 measured at time t - k relative to publications at time t. We lag the instrument because there will be some delay in receiving research funds and publishing a research paper – individual lags are included up to k=10.  $\psi^h(L)$  is a polynomial of the lag operator, and  $z_{t-1}$  is a vector of control variables, which includes the same controls included in the second

#### Figure 5: Relationship between National Funding and National Mortality



Notes: This figure shows the relationship between the share of national funding allocated to a disease at time t and three leads and lags of the national share of YPLL of that disease. Equation 4 shows the specification.

stage regression. The fixed effects included in 1 are also included. These are for location\*time  $(\delta_{lt})$ , which account for changing demographics in a local area; disease\*time  $(\gamma_{dt})$ , which captures national trends in the disease; location\*disease  $(\eta_{ld})$ , which capture level differences in the mortality rate for a particular disease in a particular place.

Figure 6 plots first stages estimates. Each coefficient is the relationship between publications and the funding instrument ( $\beta$ ) estimate using a different lag length k. The F-statistics associated with the regression is also shown. We select the first stage with the largest F-statistic, which is a lag length of 2.

#### 4.2.2 ARRA shock

The previous instrument is identified using fluctuations in national funding across diseases. As an alternative approach, we use as an instrument a single sharp funding shock. In 2009 and 2010, the NIH budget was increased by \$8 billion dollars as part of the American Recovery and Reinvestment ACT (ARRA). The ARRA was part of the stimulus package directed at the economic downturn resulting from the Global Financial Crisis. This funding was allocated to research through two broad avenues: "ARRA Solicited" and "Not ARRA Solicited", that are identified explicitly in the data. The former was a call for applications and the latter was essentially a retro-active extension of the payline for proposals received in the fiscal years 2008 and 2009. Here we will use the "Not ARRA Solicited" funding mechanism as a shock to local funding for a particular disease. This funding represents an unexpected funding "windfall" to the local area.<sup>13</sup>

<sup>&</sup>lt;sup>13</sup>We exclude new FOA because part of the purpose of ARRA was economic recovery and areas that were hurt more economically could have experience a worse health shock and hence the health recovery and stimlus could cause reversion to the mean.





Notes: This figure shows the estimates of Equation 6 for a range of lags and the F-statistic for each regression. The red markes gives the F-stats (right axis) and grey the elasticity of pubs to imputed availability of funding (left axis). Each coefficient is estimated using a separate regression. The lag length k indicates the timing of the instrument relative to publications.

**Exclusion Restriction** The exclusion restriction requires that the unobserved component of local mortality is not related to the increase in funding from the ARRA shock. Importantly, because the ARRA was a stimulus package and health can be affected by negative economic shocks, regression to the mean is of concern with this shock. The component of the ARRA shock that we use in the analysis is the expansion of the payline for 2008 and 2009 grant applications.

There are two identification concerns. Firstly, it is possible that lowering the payline extends funding in a systematic manner towards places with better or worse local health (depending on how this is related to the proposal score). To examine this, we relate the funding expansion to the mortality trends from the pre-period 2001-2007. The left panel of Figure 7 shows a binned scatter plot of the windfall funding gain against the six year growth rate in age-adjusted YPLL per capita, controlling for location and disease fixed effects. The relationship is negative, although not statistically significant, indicating that faster growth in mortality was on average associated with receiving less windfall funding. Secondly, since ARRA was an economic stimulus package, it is possible that it went to locations where the economic impacts of the financial crisis were worse, which could be locations that experience worse health shocks. To examine this possibility the right panel of 7 relates the percentage change in the YPLL from 2008 to 2009 with the windfall funding. The relationship is slightly negative, which means that places where YPLL was bigger (which is an adverse health change) received less funding on average through the ARRA.

These statistics suggest that if anything the concern should be that funding is going to diseaselocation combinations that are already trending towards lower mortality. To examine this possibility in greater depth, we repeat this exercise breaking down the relationship by disease (Figure 8). We find that the negative relationship does not hold across diseases - for some diseases funding is allocated to HRRs that were trending towards improvement, while in others the trend was towards worse health and some have no relationship. This suggests that the NIH is not systematically

#### Figure 7: Relationship between ARRA windfall and Growth in YPLL



Notes: This figure shows a bin scatter of the windfall funding received as a result of the ARRA shock (relative to regular funding in 2009) against the growth rate of YPLL from 2008-2009 and from 2001-2007.

allocating marginal funding on the basis of local health trends.

**First stage** The first stage regression is as follows:

$$log(R_{dl,t}) = \alpha + \beta log(ARRAfunding)_{dl,t-k} + \psi(L)z_{t-1} + \delta_{lt} + \gamma_{dt} + \eta_{ld} + \epsilon_{ld,t}$$
(6)

We pool together ARRA funding received in 2009 and 2010 so as not to identify the effects based on differences in the disbursal of funding across these two years. This variable is equal to zero in every year except 2009. We control for three lags of non-ARRA funding. Figure 9 shows the estimates of the  $\beta$  and corresponding F-statistics for each regression. The F-statistic is highest for two lags and so we use this as our first stage.

#### 4.3 Results

We first estimate IRFs using the instruments directly in the regression (i.e the reduced form, estimated using OLS). Figure 10 shows the reduced form relationship between the mortality outcome and the instruments (full regression results are shown in Tables A.3 and A.7). The funding shock is lagged three periods relative to the outcome variable. This means that the coefficients measure the effect of a 1% shock to research funding three years ago. Unlike the OLS relationship between age-adjusted years of potential life lost and both publications and NIH funding, the relationship between the funding shocks and and age-adjusted years of potential life lost is clearly negative. The shape of the IRF is in line with our hypothesis that research findings are adopted first in the area where they were produced but eventually spread geographically. Since we estimate the mortality effect relative to other geographic locations, as the research results spread the relative benefit to the location where they were first produced will disappear.



Figure 8: Relationship between ARRA windfall funding and pre-period mortality changes by disease

Notes: This figure shows a bin scatter of the estimated elasticity of the windfall funding received as a result of the ARRA shock (relative to regular funding in 2009) against the growth rate of YPLL from 2001-2007.





Notes: This figure shows the estimates of Equation 6 for a range of lags and the F-statistic for each regression. The red markes gives the F-stats (right axis) and grey the elasticity of pubs to imputed availability of funding (left axis). Each coefficient is estimated using a separate regression. The lag length k indicates the timing of the instrument relative to publications.

Figure 10: Reduced form estimates of impulse response function of the effect of funding shock instruments directly on years of potential life lost



(a) Funding Capture Instrument

Notes: These figures show the IRF from an OLS estimates of Equation 1 using the relevant instrument directly as the measure of the shock. Panel (a) shows the results for the instrument that uses the local elasticity of funding as weight to national shocks in disease funding. Panel (b) uses the ARRA extension of the payline as the instrument for research publications. In both cases the instrument is a funding shock that occured two years prior to the publication output, and three years prior to the mortality outcome. The standard errors are clustered by disease-hrr.

The IRFs from estimating Equation 1 using 2SLS are depicted for each of the instruments in Figure 11 (the full regression results are shown in Tables A.4 and A.5). Using the funding capture instrument (Figure 11 (a)), a one percent shock to publications on a specific disease produced by local researchers reduces the mortality rate for that disease by a cumulative amount of 0.35%. This occurs across the first five years following the publication of the research. At the mean, a one percent increase in publications represents an increase of 5.8 publications and the mean years of potential life lost is 6,450. Hence at the mean one extra publication reduces years of potential life lost by 388 years.

Using the ARRA shock (Figure 11 (b)) produces a negative effect of 0.06% in the first year following the shock. The effect fades more quickly than using the funding capture instrument. Unlike the funding capture instrument, there is a second statistically significant reduction in mortality 5 years following the shock.

A question of policy interest is the return to investmet on research funding. We estimate 1 using 2SLS with NIH funding as the endogenous variable and the funding capture instrument. The results are shown in Figure 12. Cumulatively, a one percent increase in local funding reduces local mortality from that disease by 0.22%. This is similar in magniture to the estimate we obtained using publications as the endogenous variable. The shape of the impulse response is broadly similar, which is a artifact of using the same instrument. However, because funding affects publications with a lag the effects continue further into the past. In Figure 11, the funding shock occurs three years prior to the outcome. Here, it occurs only one year prior because we measure the shock in the same period as the endogenous measure of funding. Since we use the same instrument for both, we are not able to truly distinguish between the effect of research outputs and inputs on local health.

There are two ways we can compute the return to research funding. The first is using the direct estimates of local funding on mortality. Here, at the mean of funding and years of potential life lost, a \$100,000 increase in funding on average saves 401 life years. Alternatively, we can take the average NIH funding per funded publication and use our estimate of the return to a publication on years of potential life. A funded publications on average recieves \$131,868 in funding. The return on a \$100,000 investment in research in a local area hence produces 294 years of savings in life. There are several important caveats to note here. First, this estimate does not include spillovers of NIH funding onto non-funded research. Second, it assumes that NIH research is producing the same health benefit on average as other publications, however, Figure 3 suggests it might actually be the more impactful research.

# 5 Heterogeniety

Among our thirty-eight diseases, 19 are cancers. It is possible that our results are driven by advances in particular diseases and previous research by Lichtenberg (2018a) has found that cancer survival has improved due to research. In this section we investigate the role that cancer plays

Figure 11: IRF: Second stage estimates of impulse response function of the effect of publications on years of potential life lost



(a) Funding Capture Instrument

Notes: These figures show the IRF from a TSLS estimates of Equation 1 of YPLL on publications. Panel (a) uses the instrument constructed using the local elasticity of funding to national funding to identify shocks in disease funding. Panel (b) uses the ARRA extension of the payline as the instrument for research publications. The standard errors are clustered by disease-hrr.

Figure 12: IRF: Second stage estimates of impulse response function of the effect of NIH funding on years of potential life lost



Notes: These figures show the IRF from a TSLS estimates of Equation 1 of YPLL on NIH funding. NIH funding is instrumented using the funding capture instrument. The instrument is measured in the same time period as the endogenous variable. The standard errors are clustered by disease-hrr.

in generating the results. We do so by estimating the OLS, first stage and IV (using the funding capture instrument) by re-estimating Equation 1 for two groups: cancer and other diseases. We do this by interacting a group dummy with the research output and instrumenting using the interaction of this dummy with the funding capture instrument, Equation 7 shows the OLS specification.

$$log(m_{ld,t+h}) = \alpha^{h} + \beta_{1}^{h} log(R_{dlt}) \times cancer_{d} + \beta_{2}^{h} log(R_{dlt}) + \psi^{h}(L)z_{t-1} + \delta_{lt}^{h} + \gamma_{dt}^{h} + \eta_{ld}^{h} + \epsilon_{ld,t+h}, for \ h = (0, 1, 2, ...10)$$
(7)

Figure 13 shows the IRF for cancer  $(\hat{\beta}_1^h + \hat{\beta}_2)$  and for other diseases  $(\hat{\beta}_2)$ . The first column shows the OLS estimates of a one percent increase in local publications for cancer/other diseases on the local mortality for that group of disease. It appears that the relationship between publications and mortality is positive for cancer and negative for the group of other diseases. Once we instrument using the funding capture instrument the relationship for both groups is negative. Unlike in the main results, for cancer, it is several years before a shock to publications has an effect on reducing mortality while for the rest of the diseases the effect is immediate and then fades, similar to the main result. There are two ways to interpret the shape of the cancer IRF compared to the main result, or the result for the result of the diseases grouped together, in the context of our framework.





Notes: These IRF are estimated using Equation 7. The response for the cancer group is  $\hat{\beta}_1^h + \hat{\beta}_2^h$ , other is $\hat{\beta}_2^h$ .

The first is that it take longer for new research to be adopted into clinical practice for cancer related relative to other diseases. Second, it could instead (or also) take longer for the changes made to clinical practice as a result of new research to affect mortality. This could be seen, for example, if earlier interventions improved, the results of which would not be seen potentially until years later. Given that our measure of research is broad, there is not enough information to draw a conclusion about the specific explanation for the cancer IRF. We leave the investigation into this difference for future research.

# 6 Conclusion

This paper provides evidence that there are local spillovers from biomedical research onto local mortality. The presence of such spillovers indicates that there are unrealized gains from scientific discovery that if disseminated could improve health in the national population. In the past decade there has been recognition that translation of science in medicine is important. Our results suggest that in addition to translating basic science into treatments, there should be focus on disseminating valuable research findings (which may or may not have commercial value).

There are several limitations that are important to note. Our approach calculates a regional return on investment for NIH spending on research. Translating our estimate of the regional return on investment into an aggregate return on investment or a measure of welfare is challenging in our framework. In the former case, converting a cross-sectional multiplier into an aggregate multiplier is difficult without a general equilibrium model. Many locations are receiving shocks at the same time - if a shock from one region was taken to another region would it be cumulative or would one substitute the other? In the latter case, welfare is challenging to capture because while we could compare some imperfect measure of cumulative benefits to the cost of a publication, we are not able to incorporate other necessary elements of the welfare calculation, most notably the cost of treatment. A limitation of our dataset is that we are not able to demonstrate the mechanism by which local research spillovers onto health nor if it is truely the research output or the input of local funding that improves health. Moreover, there are other dimensions to health beyond mortality, which should be investigated in future research.

# References

- Agha, L. and Molitor, D. (2018) The local influence of pioneer investigators on technology adoption: Evidence from new cancer drugs, *The Review of Economics and Statistics*, **100**, 29–44.
- Azoulay, P., Graff Zivin, J. S. and Wang, J. (2010) Superstar extinction, The Quarterly Journal of Economics, 125, 549–589.
- Azoulay, P., Zivin, J. S. G., Li, D. and Sampat, B. N. (2018) Public r&d investments and private-sector patenting: Evidence from NIH funding rules, *The Review of Economic Studies*, 86, 117–152.
- Baicker, K. and Chandra, A. (2010) Understanding agglomerations in health care, in Agglomeration economics, University of Chicago Press, pp. 211–236.
- Bania, N., Eberts, R. W. and Fogarty, M. S. (1993) Universities and the startup of new companies: can we generalize from route 128 and silicon valley?, *The review of economics and statistics*, pp. 761–766.
- Beeson, P. and Montgomery, E. (1993) The effects of colleges and universities on local labor markets., *Review of Economics & Statistics*, 75, 753–761.
- Bloom, N., Jones, C. I., Van Reenen, J. and Webb, M. (2020) Are ideas getting harder to find?, American Economic Review, 110, 1104–1144.
- Borjas, G. J. and Doran, K. B. (2012) The collapse of the soviet union and the productivity of american mathematicians, *The Quarterly Journal of Economics*, **127**, 1143–1203.
- Borjas, G. J. and Doran, K. B. (2015) Which peers matter? the relative impacts of collaborators, colleagues, and competitors, *Review of economics and statistics*, 97, 1104–1117.
- Ching, A. and Ishihara, M. (2010) The effects of detailing on prescribing decisions under quality uncertainty, *QME*, **8**, 123–165.
- Ching, A. T., Clark, R., Horstmann, I. and Lim, H. (2016) The effects of publicity on demand: The case of anticholesterol drugs, *Marketing Science*, 35, 158–181.
- Chintagunta, P. K., Goettler, R. L. and Kim, M. (2012) New drug diffusion when forward-looking physicians learn from patient feedback and detailing, *Journal of Marketing Research*, 49, 807–821.
- Coleman, J., Katz, E. and Menzel, H. (1957) The diffusion of an innovation among physicians, Sociometry, 20, 253–270.

- Feldman, M. P. and Kogler, D. F. (2010) Chapter 8 stylized facts in the geography of innovation, in Handbook of The Economics of Innovation, Vol. 1 (Eds.) B. H. Hall and N. Rosenberg, North-Holland, vol. 1 of Handbook of the Economics of Innovation, pp. 381 – 410.
- Glaeser, E. L., Kallal, H. D., Scheinkman, J. A. and Shleifer, A. (1992) Growth in cities, *Journal of political economy*, 100, 1126–1152.
- Ham, J. C. and Weinberg, B. A. (2021) Novelty, knowledge spillovers and innovation: Evidence from nobel laureates, Tech. rep., Working Paper.
- Jordà, Ò. (2005) Estimation and inference of impulse responses by local projections, *American Economic Review*, **95**, 161–182.
- Kantor, S. and Whalley, A. (2009) Do universities generate agglomeration spillovers? evidence from endowment value shocks, Tech. rep., National Bureau of Economic Research.
- Kibria, A., Mancher, M., McCoy, M. A., Graham, R. P., Garber, A. M., Newhouse, J. P. et al. (2013) Variation in health care spending: target decision making, not geography.
- Krugman, P. (1991) Increasing returns and economic geography, Journal of political economy, 99, 483-499.
- Lam, A. (2011) What motivates academic scientists to engage in research commercialization: 'gold', 'ribbon' or 'puzzle'?, *Research Policy*, **40**, 1354–1368.
- Lichtenberg, F. R. (2013) The effect of pharmaceutical innovation on longevity: patient level evidence from the 1996–2002 medical expenditure panel survey and linked mortality public-use files, in *Forum for health economics & policy*, De Gruyter, vol. 16, pp. 1–33.
- Lichtenberg, F. R. (2018a) The impact of biomedical research on us cancer mortality, Measuring and Modeling Health Care Costs, 76, 475.
- Lichtenberg, F. R. (2018b) The impact of new drug launches on life-years lost in 2015 from 19 types of cancer in 36 countries, Journal of Demographic Economics, 84, 309–354.
- Lichtenberg, F. R. (2019a) How many life-years have new drugs saved? a three-way fixed-effects analysis of 66 diseases in 27 countries, 2000–2013, International health, 11, 403–416.
- Lichtenberg, F. R. (2019b) The long-run impact of new medical ideas on cancer survival and mortality, *Economics of Innovation and New Technology*, 28, 722–740.
- Lucas Jr, R. E. (1988) On the mechanics of economic development, Journal of monetary economics, 22, 3-42.
- Merton, R. K. (1968) The matthew effect in science: The reward and communication systems of science are considered, *Science*, **159**, 56–63.
- Moretti, E. (2004) Estimating the social return to higher education: evidence from longitudinal and repeated crosssectional data, *Journal of econometrics*, **121**, 175–212.
- Murphy, K. M. and Topel, R. H. (2006) The value of health and longevity, Journal of Political Economy, 114, 871–904.
- Myers, K. (2020) The elasticity of science, American Economic Journal: Applied Economics, 12, 103–34.
- Packalen, M. and Bhattacharya, J. (2017) Neophilia ranking of scientific journals, Scientometrics, 110, 43-64.
- Rogers, E. M. (1962) Diffusion of innovations, New York, Free Press of Glencoe.

- Romer, P. M. (1986) Increasing returns and long-run growth, Journal of political economy, 94, 1002–1037.
- Sampat, B. N. and Lichtenberg, F. R. (2011) What are the respective roles of the public and private sectors in pharmaceutical innovation?, *Health Affairs*, **30**, 332–339.
- Sattari, R. and Weinberg, B. A. (2017) Public research funding and scientific productivity.
- Saxenian, A. (1996) Regional Advantage: Culture and Competition in Silicon Valley and Route 128, With a New Preface by the Author, Harvard University Press.
- Stephan, P. E. (2010) The economics of science, in *Handbook of the Economics of Innovation*, Elsevier, vol. 1, pp. 217–273.
- Toole, A. A. (2007) Does public scientific research complement private investment in research and development in the pharmaceutical industry?, *The Journal of Law and Economics*, **50**, 81–104.
- United States Mortality DataBase (1999-2017) Available at usa.mortality.org (data downloaded on February 10 2020).
- Waldinger, F. (2010) Quality matters: The expulsion of professors and the consequences for phd student outcomes in nazi germany, *Journal of political economy*, **118**, 787–831.
- Waldinger, F. (2012) Peer effects in science: Evidence from the dismissal of scientists in nazi germany, The Review of Economic Studies, 79, 838–861.
- Zolas, N., Goldschlag, N., Jarmin, R., Stephan, P., Owen-Smith, J., Rosen, R. F., Allen, B. M., Weinberg, B. A. and Lane, J. I. (2015) Wrapping it up in a person: Examining employment and earnings outcomes for ph. d. recipients, *Science*, **350**, 1367–1371.
- Zucker, L. G., Darby, M. R. and Brewer, M. B. (1998) Intellectual human capital and the birth of us biotechnology enterprises, *The American Economic Review*, 88, 290–306.

# **Online Appendix**

# A Data Details

# A.1 Mortality

Administrative death records were obtained from the National Center for Health Statistics Mortality Multiple Cause Files 1985–2018. This dataset contains information from death certificates on up to 20 contributing factors in a death as well as the physicians opinion on the primary cause of death. Causes of death are classified using ICD-10 codes for the period 1999-2017

We use causes of death that can be matched to mesh codes. ICD10 codes are mapped to mesh codes using the Mesh on Demand generator as well as manual verification. Table A.1 shows the cross-walk that was used.

NCHS-113	Disease Group	ICD10	Mesh codes
GR113-003	Certain other intestinal infec-	A04,A07-A09	D007410, D019350, D005759, D054324, D054307, D005873, D004761,
	tions		D021866, D001447, D009663, D016360, D003457, D003092, D007411,
			D019453, D004751, D021865, D015008, D004927, D000069981,
			D054308, D012400, D002167
GR113-010	Septicemia	A40-A41	D016470, D018805, D012772
GR113-015	Viral hepatitis	B15-B19	D019698, D006505, D003699, D006521, D016751, D006506, D006509,
			D019694, D019701, D006525, D006526
GR113-016	Human immunodeficiency	B20-B24	D017088, D039682, D016263, D000386, D019053, D015658, D015526,
	virus (HIV) disease		D000163, D019247, D020943, D006679, D000071297
GR113-020	Malignant neoplasms of lip,	C00-C14	D007972, D010157, D005887, D008048, D013365, D010610, D009062,
	oral cavity and pharynx		D014062, D010307, D007012, D009303, D012468, D009959, D013362,
			D014067
GR113-021	Malignant neoplasm of esoph-	C15	D004938
	agus		
GR113-022	Malignant neoplasm of stom-	C16	D013274
	ach		
GR113-023	Malignant neoplasms of colon,	C18-C21	D005736, D001005, D003123, D012004, D003110, D015179, D011125,
	rectum and anus		D012811
GR113-024	Malignant neoplasms of liver	C22	D018281, D018197, D018285, D008113
	and intrahepatic bile ducts		
GR113-025	Malignant neoplasm of pan-	C25	D010190
	creas		
GR113-027	Malignant neoplasms of tra-	C33-C34	D001984, D008175, D010178
	chea, bronchus and lung		
GR113-028	Malignant melanoma of skin	C43	D018328, D018327, D008545

GR113-029	Malignant neoplasm of breast	C50	D058922, D018567, D064726, D001943, D000069584
GR113-030	Malignant neoplasm of cervix uteri	C53	D002583
GR113-031	Malignant neoplasms of cor- pus uteri and uterus, part un- specified	C54-C55	D016889, D014594
GR113-032	Malignant neoplasm of ovary	C56	D010051
GR113-033	Malignant neoplasm of prostate	C61	D064129, D011471
GR113-034	Malignant neoplasms of kid- ney and renal pelvis	C64-C65	D030321, D017624, D007680
GR113-035	Malignant neoplasm of blad- der	C67	D001749
GR113-036	Malignant neoplasms of meninges, brain and other parts of central nervous system	C70-C72	D018316, D019574, D020339, D005910, D003390, D001932, D016545, D020295, D013120, D005909, D002528, D007029, D004806, D009837, D016543, D001254, D002551
GR113-039	Non-Hodgkin lymphoma	C82-C85	D016411,D016403,D008224,D020522,D016393,D017728,D000069293,D018442,D016400,D016483,D002051,D008228,D054685,D016399
GR113-040	Leukemia	C91-C95	D004915, D007947, D015466, D015463, D001752, D007938, D007945, D015472, D007946, D007948, D015458, D015471, D015465, D015461, D015459, D007943, D015470, D023981, D007952, D015464, D015448, D015456, D015473, D015452, D015477, D007951, D015479, D054438, D054403, D054429, D015451
GR113-041	Multiple myeloma and im- munoproliferative neoplasms	C88,C90	D009101, D008258, D007161, D006362, D007952

GR113-045	Anemias	D50-D64	D017086, D017085, D012805, D018798, D000751, D012010, D000750,
			D000753, D000756, D000741, D000745, D000747, D000749, D000740,
			D000748, D000746, D000754, D013789, D000743, D006445, D000744,
			D005330, D005236, D005199, D006450, D000752, D006463, D004612,
			D000755, D013103, D005331, D006457, D005955, D000742
GR113-046	Diabetes mellitus	E10-E14	D003925, D014929, D003923, D003924, D003926, D003922, D003928,
			D011236, D016883, D003930, D003920, D006944, D048909, D003929,
			D017719, D058065, D000071698
GR113-047	Nutritional deficiencies	E40-E64	D002796, D013540, D052879, D010383, D001602, D007732, D020138,
			D011488, D014808, D013231, D026681, D001206, D000067011,
			D048070, D003677, D001361, D000752, D014806, D014813, D044342,
			D013217, D014804, D014899, D011191, D010018, D014811, D005494,
			D053098, D006475, D013832, D014802, D008275
GR113-050	Meningitis	G00,G03	D008590, D001100, D008216, D008584, D020814, D008582, D016921, D008586, D016919, D008586, D008586, D016919, D008586, D016919, D008586, D016919, D008586, D016919, D008586, D016919, D008586, D016919, D008586, D008586, D016919, D008586, D008566, D0085666, D0085666, D0085666, D0085666, D0085666, D008566, D008566, D0085
GR113-051	Parkinson disease	G20-G21	D010300
GR113-052	Alzheimer disease	G30	D000544
GR113-059	Acute myocardial infarction	I21-I22	D009203, D000072658, D056989, D000072657, D056988
GR113-064	Other heart diseases I26-I51	D054143,D004698,I	D009205, D010493, D010494, D006333, D059905, D004696, D004697, D054144
GR113-070	Cerebrovascular diseases	I60-I69	D020521, D009072, D002532, D020225, D020293, D002542, D046589,
			D020300, D000074042, D020766, D002539, D013345, D020299,
			$D016893,\ D055955,\ D020526,\ D002537,\ D020765,\ D002341,\ D014854,$
			D020200, D046648, D002544, D020144, D028243, D020301, D059345,
			D006408, D020146, D059409, D002543, D020767, D020145, D002545,
			D020243, D020244, D020199, D016657, D020762, D020520, D002561,
			D006407
GR113-071	Atherosclerosis	I70	D050197

GR113-082	Chronic lower respiratory dis-	J40-J47	D004646
	eases		
GR113-090	Peptic ulcer	K25-K28	D013276, D010437, D004381, D010439
GR113-094	Alcoholic liver disease	K70	D008108, D008104, D005235, D006519
GR113-097	Nephritis, nephrotic syn-	N00-N07,N17-	D009394, D015433, D009393, D015432, D007676, D005921, D005923
	drome and nephrosis	N19,N25-N27	
GR113-135	Complications of medical and	Y40-Y84,Y88	D011183
	surgical care		

Years of potential life lost are computed using the difference between life expectancy at age of death.<sup>14</sup> This years of potential life lost are summed by disease and HRR region. As is customary we assume years lived in excess of life expectancy represent 0.

$$YPLL = LE - age$$

#### A.2 Publications

The 2018 annual baseline files of the PubMed database are used to obtain information on publications that related to the diseases in Table A.1 shows the cross-walk that was used. Publications are matched using these mesh codes. A publication can count towards more than one disease if it is tagged with mesh codes that match both diseases. Publications are geolocated based on the affiliation of the first author in Author-ity! 2018.

#### A.3 Grants

Data on NIH grants is obtained from the NIH ExPORTER database. The RePORTER Project Data from FY 1985-2017, RePORTER Publications released in calendar years 1985-2017 and RePORTER Publications link tables are used to assign grant funding to disease topics. First, grants are linked to their respective publications. These publications are then merged with PubMed records to obtain the mesh codes associated with the publications.

# **B** Regression Result Tables

This section provides the full regression results from the estimation of Equation 1 using OLS and TSLS. Table A.2 shows the results for a regression of the inverse sine of age-adjusted years of potential life lost on the inverse sine of "all publications". The top left panel of Figure 3 plots the coefficient on the first lag of publications from each regression.

<sup>&</sup>lt;sup>14</sup>https://usa.mortality.org/national.php?national=USA

	0	1	2	3	4	5	6	7	8	9
L.log(publications)	-0.001	0.000	-0.000	0.001	0.000	0.001	0.000	$0.003^{**}$	0.002	$0.004^{**}$
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
	0.000	0.000	0.001	0.000	0.001	0.000	0.000*	0.001	0.000**	0.000*
L2.log(publications)	0.000	-0.000	0.001	0.000	0.001	0.000	0.003*	0.001	0.003**	0.003*
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
L3 log(nublications)	-0.000	0.001	0.000	0.001	0.001	0.003*	0.001	0.003*	0.002*	0.002
10.10g(publications)	(0.000)	(0.001)	(0.000)	(0.001)	(0.001)	(0.000)	(0.001)	(0.000)	(0.002)	(0.002)
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
L4.log(publications)	0.001	0.000	0.000	0.001	0.002	0.000	0.002	0.002	0.001	0.001
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
	· · · ·	× ,	· · · ·	× ,	· · · ·	× ,				
$L5.\log(publications)$	0.000	0.000	0.001	0.002	0.000	0.001	0.001	0.001	0.002	-0.001
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
L6.log(publications)	0.000	0.001	0.001	-0.000	0.002	0.001	0.001	0.002	-0.001	-0.001
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
I 7 log(publications)	0.001	0.001	0.000	0.001	0.000	0.000	0.001	0.001	0.001	0.002*
L1.log(publications)	(0.001)	(0.001)	-0.000	(0.001)	(0.000)	(0.000)	(0.001)	-0.001	-0.001	-0.003
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
L8.log(publications)	0.001	-0.000	0.002	0.000	0.000	0.001	-0.001	-0.001	-0.002	-0.001
)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
	(0001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
L9.log(publications)	-0.000	0.001	-0.000	0.000	-0.000	-0.001	-0.002	$-0.002^{*}$	-0.002	-0.001
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
$L10.\log(publications)$	0.001	-0.000	-0.000	-0.000	-0.001	-0.002*	-0.002	-0.002	-0.001	0.002
	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)	(0.001)
()tt	2 2006***	9 100***	9 100***	o 101***	0 170***	0 171***	9 100***	9 160***	9 1F0***	ዓ 1 <b>୮</b> / * * *
Constant	3.200	3.198	$3.190^{-1}$	3.181	3.170	3.1(1)	3.108	3.102	3.138	3.134
	(0.004)	(0.004)	(0.004)	(0.004)	(0.004)	(0.004)	(0.004)	(0.004)	(0.005)	(0.005)
Observations	2.21e+05	2.09e + 05	1.98e + 05	1.86e + 05	1.74e + 05	1.63e + 05	1.51e + 05	1.40e + 05	1.28e + 05	1.16e + 05
Fixed Effects	77	17	37	77	77	17	77	37	77	77
Year-HKK	X	X	X	X	X	X	X	X	X	X
Year-Disease	X	X	X	X	X	X	X	X	X	X
HRR-Disease	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х

Table A.2: OLS results

	0	1	2	3	4	5	6	7	8	9
L3.Instrument	-0.0123**	-0.0135**	-0.0125**	-0.00935	-0.00759	-0.00588	-0.00446	-0.00275	-0.00187	0.00169
	(0.00448)	(0.00444)	(0.00473)	(0.00486)	(0.00538)	(0.00578)	(0.00615)	(0.00616)	(0.00612)	(0.00659)
Constant	3.323***	3.327***	3.310***	$3.274^{***}$	3.252***	3.231***	3.213***	3.193***	3.180***	3.143***
	(0.0417)	(0.0413)	(0.0440)	(0.0452)	(0.0501)	(0.0538)	(0.0571)	(0.0572)	(0.0568)	(0.0612)
Observations	220894	209304	197676	186048	174420	162792	151164	139536	127908	116280
Fixed Effects										
Year-HRR	Х	Х	Х	Х	Х	Х	X	Х	Х	Х
Year-Disease	X	Х	X	X	X	X	X	X	Х	Х
HRR-Disease	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х

Table A.3: Reduced Form: Funding Capture Instrument

This table shows the results from estimating Equation 1 using OLS. Each column has the outcome for the h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the dependent variable of interest is the funding capture instrument, lagged three periods. Regressions control for all three pairs of fixed effects . Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

	0	1	2	3	4	5	6	7	8	9
L.log(publications)	-0.0916**	-0.102**	-0.0929**	-0.0689*	-0.0563	-0.0454	-0.0365	-0.0237	-0.0182	0.00923
	(0.0335)	(0.0335)	(0.0347)	(0.0342)	(0.0374)	(0.0411)	(0.0433)	(0.0429)	(0.0424)	(0.0485)
L2.log(publications)	$0.0100^{**}$ (0.00364)	$0.0106^{**}$ (0.00352)	$0.0102^{**}$ (0.00344)	$0.00708^{*}$ (0.00320)	0.00618 (0.00331)	0.00407 (0.00323)	0.00521 (0.00293)	0.00249 (0.00252)	$0.00409^{*}$ (0.00187)	0.00270 (0.00139)
L3.log(publications)	$0.00659^{*}$ (0.00258)	$\begin{array}{c} 0.00823^{***} \\ (0.00245) \end{array}$	$0.00688^{**}$ (0.00241)	$0.00560^{*}$ (0.00228)	0.00423 (0.00225)	$0.00505^{*}$ (0.00226)	0.00244 (0.00201)	$0.00363^{*}$ (0.00157)	$0.00277^{*}$ (0.00124)	$0.00208 \\ (0.00126)$
L4.log(publications)	$0.00485^{**}$ (0.00165)	$0.00429^{**}$ (0.00154)	$0.00392^{**}$ (0.00151)	$0.00285^{*}$ (0.00137)	$0.00341^{*}$ (0.00132)	0.00110 (0.00124)	0.00200 (0.00108)	$0.00169 \\ (0.00116)$	0.000986 (0.00150)	0.00172 (0.00239)
$L5.\log(publications)$	$0.00308^{*}$ (0.00129)	$0.00322^{*}$ (0.00128)	$0.00348^{**}$ (0.00124)	$0.00324^{**}$ (0.00118)	$0.00105 \\ (0.00113)$	$0.00186 \\ (0.00106)$	0.00117 (0.00108)	0.000612 (0.00127)	$0.00132 \\ (0.00165)$	-0.000570 (0.00227)
Observations	220894	209304	197676	186048	174420	162792	151164	139536	127908	116280
Fixed Effects										
Year-HRR	X	Х	Х	Х	Х	Х	Х	Х	X	Х
Year-Disease	X	Х	Х	Х	Х	Х	Х	Х	X	Х
HRR-Disease	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х

Table A.4: Instrumental Variables: Funding Capture Instrument

This table shows the results from estimating Equation 1 using TSLS and the funding capture instrument. Each column has the outcome for the h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the endogenous dependent variable of interest is on lag of the inverse sine of publications in a location (hospital market) related to each disease in each year. Regressions control for all three pairs of fixed effects and ten lags of publications. Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

	0	1	2	3	4	5
L.log(publications)	$-0.0685^{*}$	-0.0207	0.00122	0.0274	$-0.0718^{*}$	-0.0263
	(0.0310)	(0.0303)	(0.0268)	(0.0285)	(0.0286)	(0.0258)
L lfunding orm	-0 000945*	-0.000551	0.000180	-0.000112	-0 00103**	-0 0000227
$\mathbb{E}$ $\mathbb{E}_n$ or $m$	(0.000402)	(0.000378)	(0.000100)	(0.000361)	(0.00105)	(0.000347)
	(0.000102)	(0.000010)	(0.000012)	(0.000001)	(0.000001)	(0.000011)
$L2.lfunding_n orm$	0.000842	0.000559	-0.000557	-0.000412	$0.00124^{**}$	0.000218
	(0.000512)	(0.000455)	(0.000403)	(0.000423)	(0.000408)	(0.000391)
	0.00100*	0.0001.40	0.00000.4	0.000104	0.00100*	0.000.450
L3. If unding $norm$	0.00123*	0.000148	0.000294	0.000104	0.00103*	0.000479
	(0.000491)	(0.000481)	(0.000423)	(0.000441)	(0.000451)	(0.000395)
Observations	220894	209304	197676	186048	174420	162792
Fixed Effects						
Year-HRR	Х	Х	Х	Х	Х	Х
Year-Disease	Х	Х	Х	Х	Х	Х
HRR-Disease	Х	Х	Х	Х	Х	Х

Table A.5: Instrumental Variables: ARRA Instrument

This table shows the results from estimating Equation 1 using TSLS and the ARRA instrument. Each column has the outcome for the h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the endogenous dependent variable of interest is one lag of the inverse sine of publications in a location (hospital market) related to each disease in each year. Regressions control for all three pairs of fixed effects and five lags of regular funding. Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

	0	1	2	3	4	5	6	7	8	9
L.log(funding)	-0.0407**	-0.0401**	-0.0332*	-0.0314*	-0.0391*	-0.0388*	-0.0310	-0.0265	-0.0183	0.00558
/	(0.0131)	(0.0133)	(0.0142)	(0.0151)	(0.0160)	(0.0163)	(0.0171)	(0.0177)	(0.0179)	(0.0196)
$L2.\log(funding)$	0.0262**	0.0265**	$0.0203^{*}$	$0.0194^{*}$	$0.0254^{*}$	$0.0246^{*}$	0.0190	0.0167	0.0113	-0.00339
	(0.00866)	(0.00868)	(0.00922)	(0.00970)	(0.0101)	(0.0102)	(0.0106)	(0.0107)	(0.0106)	(0.0113)
L3.log(funding)	0.000628	-0.00125*	0.0000839	0.000909	-0.000500	-0.000709	0.000736	-0.000212	-0.000504	0.00143
	(0.000570)	(0.000535)	(0.000556)	(0.000591)	(0.000613)	(0.000675)	(0.000688)	(0.000671)	(0.000747)	(0.000885)
L4.log(funding)	-0.00106*	0.0000596	0.000785	-0.000454	-0.000535	0.000544	-0.000714	-0.000975	0.000494	-0.000179
	(0.000515)	(0.000553)	(0.000567)	(0.000559)	(0.000627)	(0.000680)	(0.000699)	(0.000707)	(0.000795)	(0.000941)
L5.log(funding)	-0.00153	-0.00124	-0.00179	-0.00111	-0.00113	-0.00217*	-0.00146	-0.000584	-0.000895	0.00115
	(0.000848)	(0.000864)	(0.000915)	(0.000969)	(0.00107)	(0.00110)	(0.00113)	(0.00129)	(0.00138)	(0.00155)
Observations	220894	209304	197676	186048	174420	162792	151164	139536	127908	116280
Fixed Effects										
Year-HRR	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х
Year-Disease	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х
HRR-Disease	Х	Х	Х	Х	Х	Х	Х	Х	Х	Х

Table A.6: Instrumental Variables Results for Funding: Funding Capture Instrument

This table shows the results from estimating Equation 1 using OLS. Each column has the outcome for the h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the dependent variable of interest is NIH funding. NIH fuding is instrumed using the funding capture instrument the same period as NIH funding. Regressions control for all three pairs of fixed effects. Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

	0	1	2	3	4	5
L3.log(ARRA funding)	-0.00111	-0.000398	-0.0000321	0.000490	-0.00141*	-0.000607
	(0.000624)	(0.000621)	(0.000619)	(0.000618)	(0.000616)	(0.000615)
$L.lfunding_n orm$	-0.000166	-0.000300	0.000182	-0.000390	-0.000253	0.000276
	(0.000249)	(0.000259)	(0.000269)	(0.000279)	(0.000289)	(0.000299)
L2 lfunding_orm	-0.000174	0.000274	-0.000548	-0.0000700	0.000420	-0.0000833
12	(0.000297)	(0.000306)	(0.000315)	(0.000326)	(0.000335)	(0.000345)
	(0.000_0.)	(0.000000)	(0.000010)	(0.00020)	(0.000000)	(0.0000010)
$L3.lfunding_n orm$	0.000323	-0.000352	0.0000636	0.000466	-0.0000731	-0.000146
	(0.000301)	(0.000309)	(0.000318)	(0.000326)	(0.000335)	(0.000348)
L4.lfunding <sub><math>n</math></sub> orm	-0.000310	0.0000701	0.000388	-0.0000934	-0.0000912	0.000436
	(0.000303)	(0.000311)	(0.000318)	(0.000327)	(0.000339)	(0.000354)
L5 lfunding orm	0.000250	0.000388	0.0000232	0 000194	0 000471	0 0000848
Lo.nununig <sub>n</sub> orm	(0.000250)	(0.000360)	(0.0000202)	(0.000194)	(0.000411)	(0,00000040)
	(0.000239)	(0.000204)	(0.000212)	(0.000280)	(0.000291)	(0.000303)
Constant	$3.209^{***}$	$3.201^{***}$	$3.194^{***}$	$3.187^{***}$	$3.179^{***}$	$3.174^{***}$
	(0.00148)	(0.00153)	(0.00159)	(0.00167)	(0.00175)	(0.00183)
Observations	220894	209304	197676	186048	174420	162792
Fixed Effects						
Year-HRR	Х	Х	Х	Х	Х	Х
Year-Disease	Х	Х	Х	Х	Х	Х
HRR-Disease	Х	Х	Х	Х	Х	Х

Table A.7: Reduced form: ARRA Instrument

This table shows the results from estimating Equation 1 using OLS. Each column has the outcome for the h step-ahead. The outcome is the inverse sine of age-adjusted years of potential life lost (YPLL) per capita and the dependent variable of interest is the ARRA instrument, lagged three periods. Regressions control for all three pairs of fixed effects . Standard errors are clustered at the disease-location to account for correlation across time within cross-sectional groupings.

# C Robustness checks

## C.1 Place of occurrence

This shows the main result using the place of occurrence of death rather than the place of residence at time of death.

Figure 14: IRF: Second stage estimates of impulse response function of the effect of publications on years of potential life lost: Funding Capture Instrument



Notes: This figures show the IRF from a TSLS estimates of Equation 1 of YPLL on publications. The standard errors are clustered by disease-hrr.

#### C.2 Diseases with high concentration of research

Our main results exclude three diseases that are relative rare causes of death and have very few publications, which are highly concetrated in one HRR. Figure 15 shows the main result with these three diseases included (hence there are 41 diseases). The results are very similar to Figure 11.

Figure 15: IRF: Second stage estimates of impulse response function of the effect of publications on years of potential life lost: Funding Capture Instrument, extra diseases



Notes: This figures show the IRF from a TSLS estimates of Equation 1 of YPLL on publications. The standard errors are clustered by disease-hrr.

# C.3 No restriction on publication time

This figure shows the results when publications can occur before the grant.

Figure 16: IRF: Second stage estimates of impulse response function of the effect of publications on years of potential life lost: Funding Capture Instrument



Notes: This figures show the IRF from a TSLS estimates of Equation 1 of YPLL on publications. The standard errors are clustered by disease-hrr.

## C.4 Instrument constructed excluding own funding

We construct an version of the instrument where national funding excludes the disease and locations own funding. We do this both when estimating the the elasticities  $(\hat{\beta}_{dl})$  using Equation 2and when construction the elasticity using Equation 3. The IRF for the IV is show in Figure 17. The results are very similar to the main results in Figure 11 Figure 17: IRF: Second stage estimates of impulse response function of the effect of publications on years of potential life lost: Funding Capture Instrument excluding own funding



Notes: This figures show the IRF from a TSLS estimates of Equation 1 of YPLL on publications. The standard errors are clustered by disease-hrr.